



Beowulf - a DBMS alternative for GLAST

With thanks to Richard Fink,
Code 631



Spatial Organization

- The GLAST study report recommended narrow-field FITS files each of which represents a mosaic tile on the sky.
- Reprocessing would mean reassignment of photons to other FITS tile files and their deletion from their current location.
- Queries involving indices that were NOT spatial could be very slow, e.g. find all photons $> 20\text{GEV}$
- BUT spatial queries are fast because they are sequential.
- Update in place of a particular photon would be hard.
- Commercial DBMS systems that manage photons in the FITS files would be expensive - and the update and time-order retrieval problem still exists.



Beowulf

- Beowulf creates supercomputers out of commodity PCs for low price/high performance.
- With disks on each PC, the Beowulf has a large number of independent and parallel I/O channels.
- Setting each PC to read the photons on its disk and accept/reject them maximizes throughput.
- Storing data randomly in space and sequentially in time is an efficient I/O strategy for GLAST
- Avoiding high cost DBMS licensing is a plus

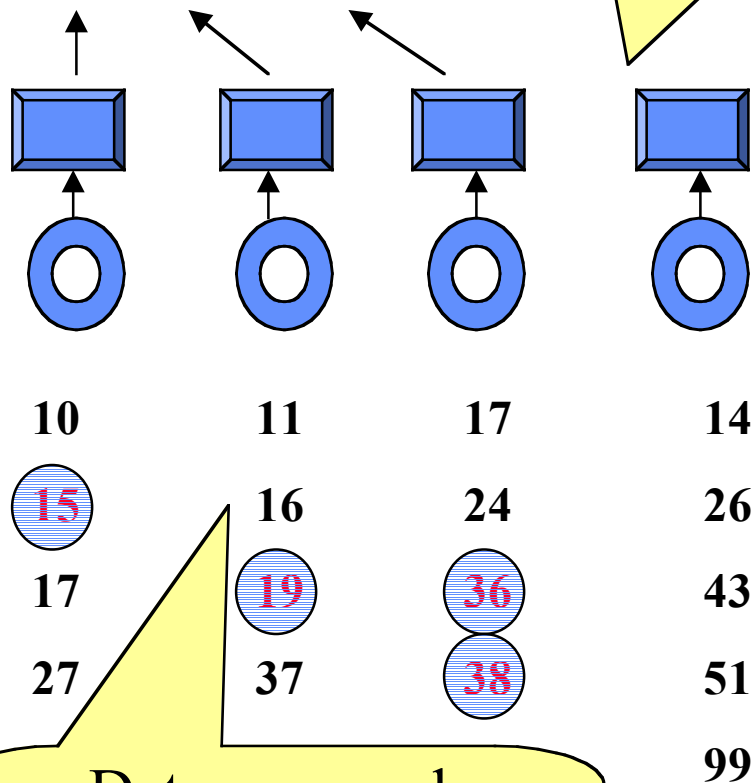
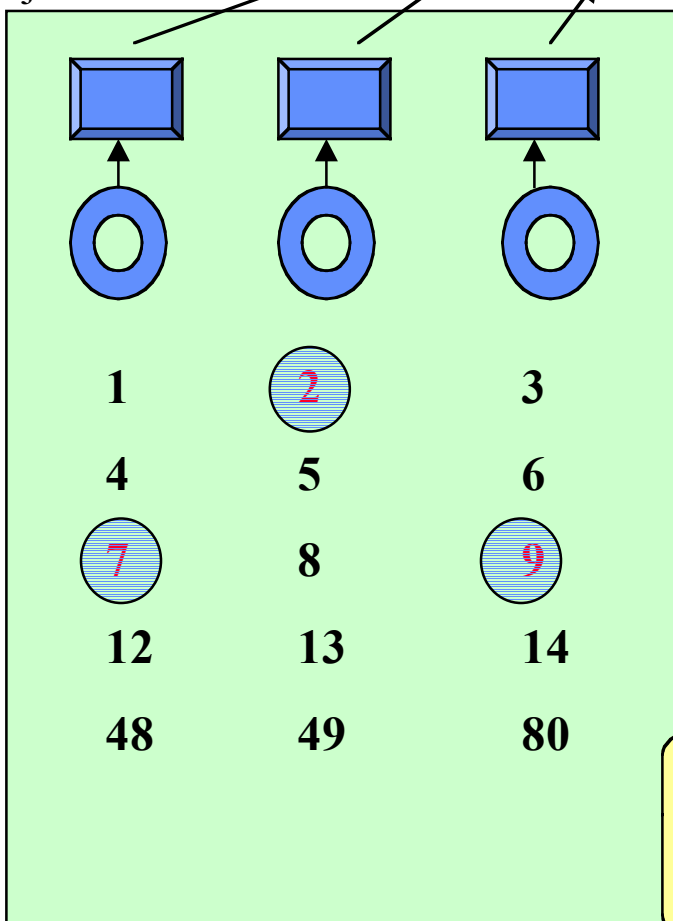


Beowulf

Thruput is max
when required
elements are
random

22 - accepted
22 - rejected

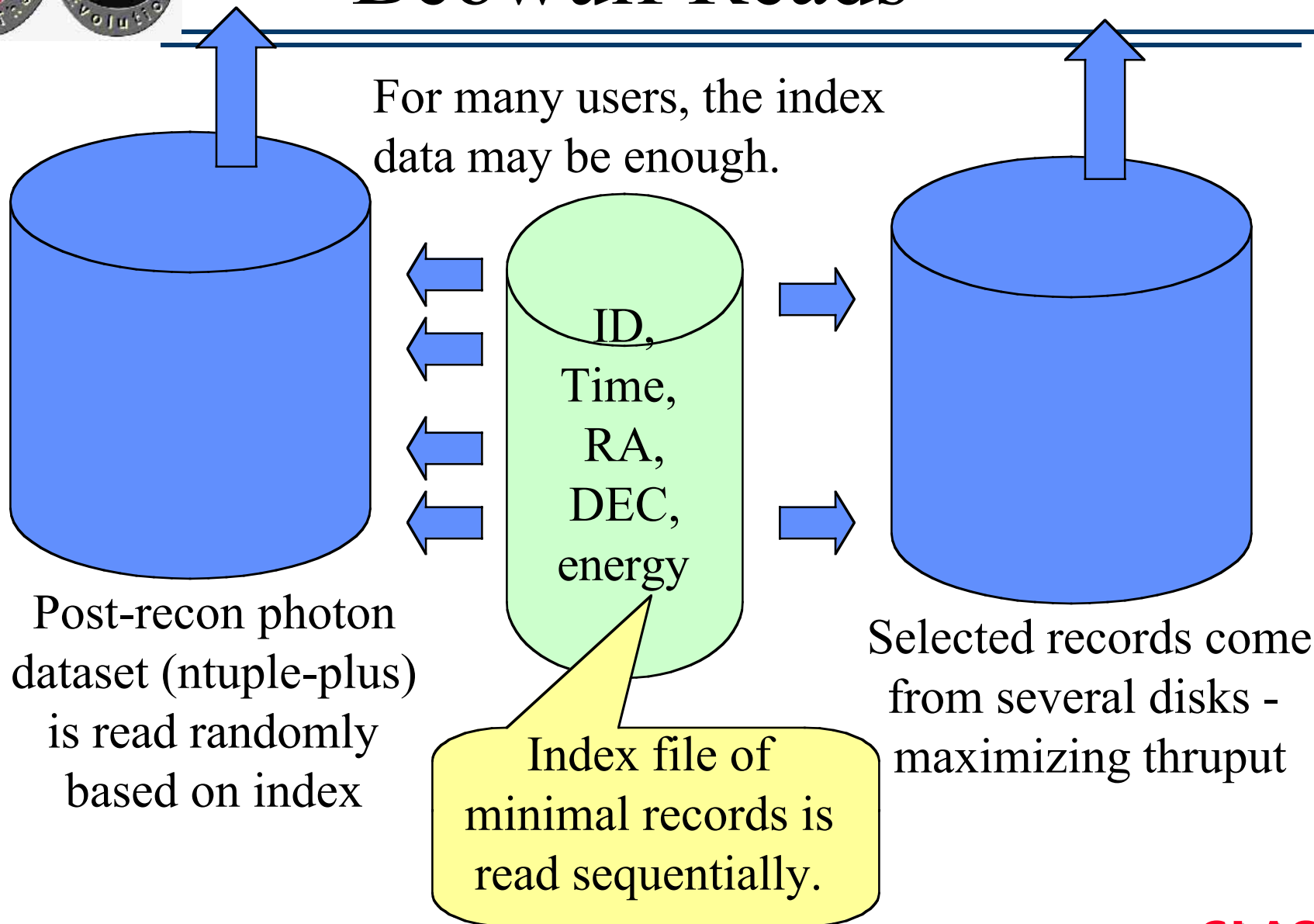
Satisfy Spatial Query



Data arranged
randomly wrt space -
sequentially in time



Beowulf Reads





The Index Table

Random reads are costly in I/O. High-yield queries will be faster read sequentially. All queries will have an initial index-read overhead. The index file for 1 year is $\sim 2\text{GB}$ - feasible for keeping in memory.

Time in secs.

<i>Events per Access</i>	<i>Index Read Time</i>	<i>Event Read Time</i>	<i>Ethernet Transfer Time</i>	<i>Total time</i>	<i>Archive Service Rate (per day)</i>
1,000	150	9	0.02	159	543
10,000	150	84	0.2	234	369
100,000	150	833	2	985	87
1,000,000	150	1680	20	1550	55
10,000,000	150	1680	200	1730	50

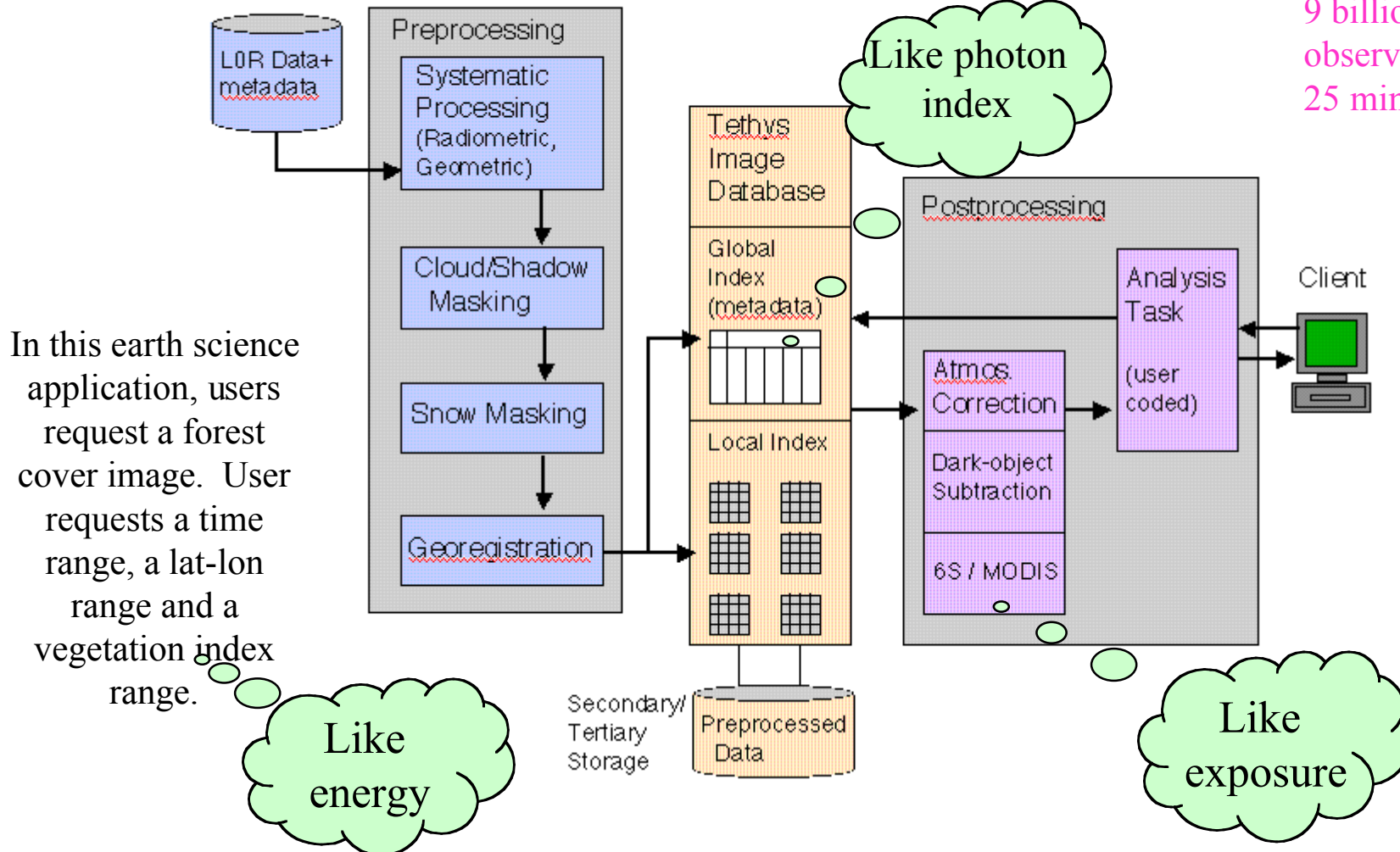
The entire index must be read, no matter how few photons are actually retrieved

Once the index is read and the yield is high, it is faster to read the data file sequentially than randomly



An example

UMD REALM Prototype



A sample query processed nearly 9 billion separate observations in 25 minutes



Upgrade Path

- Add entire computers and disks over time without changing the architecture.
- Replace components if they break
- Not have to rewrite any software
- Requires the choice of a platform and operating system that has a lot of inertia, i.e. likely to persist
- Requires an architecture that is extensible.
- Hardware replacement can be done with a distributed DBMS or Beowulf - but license cost may be a big issue with a commercial DBMS